

## 10.1 Advice for applying Machine Learning

Debugging a learning algorithm:

- Suppose you have implemented regularized linear to predict housing pricing

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (\text{h}_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^m \theta_j^2 \right]$$

However, when you test your hypothesis on a new set of houses, you find that it makes unacceptably large errors in its predictions. What should you try next?

- Get more training example  $\rightarrow$  fixes high variance
- Try smaller set of features  $\rightarrow$  fixes high variance
- Try getting additional features  $\rightarrow$  fixes high bias
- Try adding polynomial features  $\rightarrow$  fixes high bias
- Try decreasing  $\lambda$   $\rightarrow$  fixes high bias
- Try increasing  $\lambda$   $\rightarrow$  fixes high variance

## Machine Learning diagnostic.

Diagnostic: A test that you can run to gain insight what is/ isn't working with a learning algorithm, and gain guidance as to how best to improve performance

Diagnostic can take time to implement, but doing so can be very good use of time

## 10.2 Evaluating your hypothesis

Training / testing procedure for linear regression

- Learn parameter  $\theta$  from training data (minimizing training error  $J(\theta)$ )
- Compute test set error:

$$J_{\text{test}}(\theta) = -\frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} y_{\text{test}}^{(i)} \log h_{\theta}(x_{\text{test}}^{(i)}) + (1 - y_{\text{test}}^{(i)}) \log h_{\theta}(x_{\text{test}}^{(i)})$$

- Misclassification error (0/1 misclassification error):

$$\text{err}(h_{\theta}(x), y) = \begin{cases} 1 & \text{if } h_{\theta} \geq 0.5, y=0 \\ 0 & \text{if } h_{\theta} < 0.5, y=1 \end{cases}$$

$$\text{Test error} = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \text{err}(h_{\theta}(x_{\text{test}}^{(i)}), y_{\text{test}}^{(i)})$$

## 10.3 Model Selection and Train Validation

$$1. h_{\theta}(x) = \theta_0 + \theta_1 x$$

$$2. h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

$$3. h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

⋮

$$10. h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10}$$

Choose  $\theta_0 + \dots + \theta_5 x^5$

How well does the model generalise? Report test set error  
 $J_{\text{test}}(\theta^{(5)})$

Problem:  $J_{\text{test}}(\theta^{(5)})$  is likely to be an optimistic estimate of generalisation error. i.e our extra parameter ( $d = \text{degree of polynomial}$ ) is fit to test set.

## Evaluating the hypothesis

60% - training set

20% - cross validation set

20% - test set

## Train / Validation / test error

Training error:

$$J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Cross validation error:

$$J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^m (h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)})^2$$

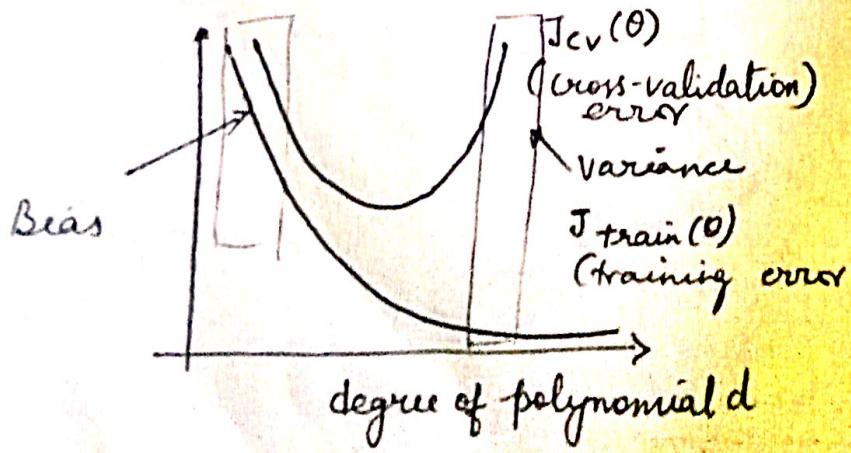
Test error:

$$J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^m (h_{\theta}(x_{\text{test}}^{(i)}) - y_{\text{test}}^{(i)})^2$$

- Select the model which give least validation error
- Estimate the generalisation error for the model

## 10.4 : Diagnosing Bias Vs Variance

Suppose your learning algorithm is performing less well than you were hoping ( $J_{\text{cv}}(\theta)$  or  $J_{\text{test}}(\theta)$  is high). Is it bias problem or a variance problem?



Bias (underfit)

- $J_{\text{train}}(\theta)$  will be high
- $J_{\text{cv}}(\theta) \approx J_{\text{train}}(\theta)$

or Variance (overfit)

- $J_{\text{train}}(\theta)$  will be low
- $J_{\text{cv}}(\theta) \gg J_{\text{train}}(\theta)$

## 10.5 Linear regression with regularization

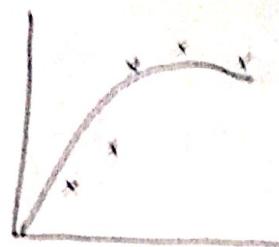
Model:  $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2$$

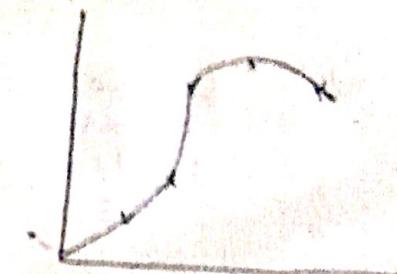


size  
Large  $\lambda$

High bias (underfit)  
 $\lambda = 1000, \theta_1 \approx 0, \theta_2 \approx 0$



size  
Intermediate ( $\lambda$ )  
Just right



size  
Small ( $\lambda$ )  
High variance (overfit)

Choosing the regularization parameter  $\lambda$

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2$$

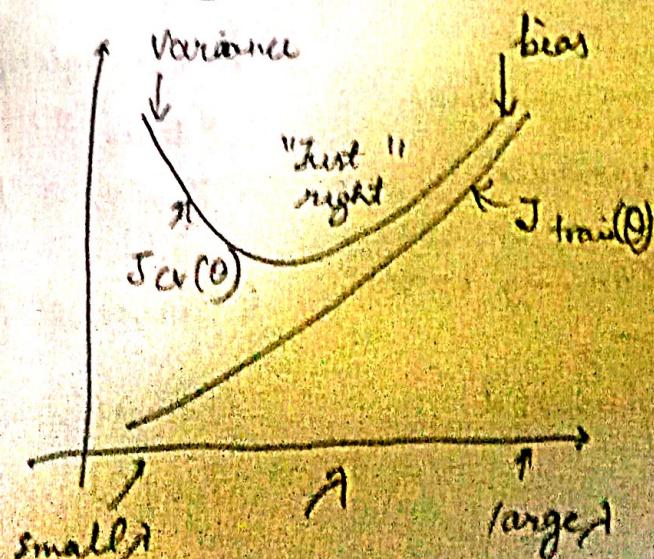
$$\text{Try } \lambda = 0 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(0)} \rightarrow J_{cv}(\theta^{(0)})$$

$$\lambda = 0.01 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$$

$$\lambda = 0.02$$

$$\lambda = 0.08 \rightarrow \dots \rightarrow \theta^{(12)} \rightarrow J_{cv}(\theta^{(12)})$$

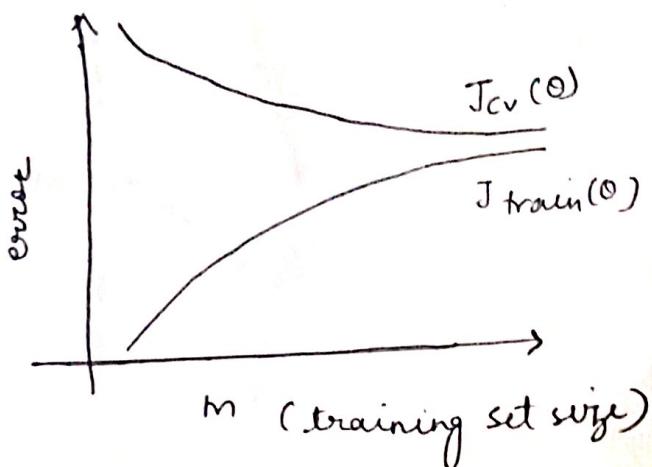
Pick which give least cross validation error



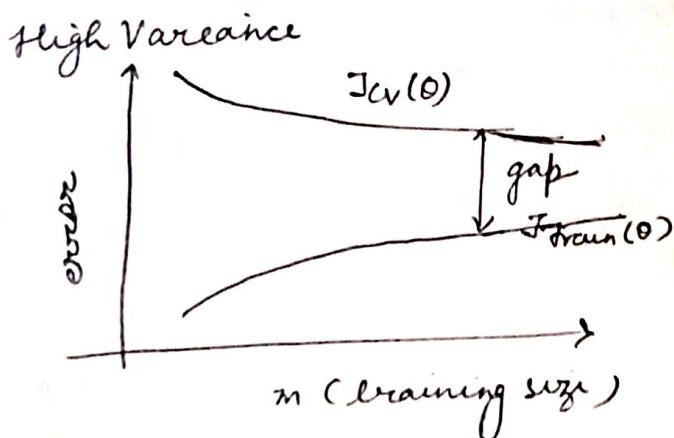
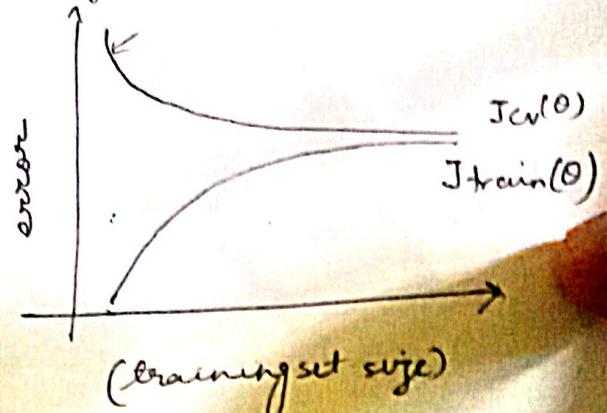
## 10.6 Learning Curves

$$J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$J_{\text{cv}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$



High bias



If a learning algorithm is suffering from high variance, getting more training data is likely to help

If learning algorithm is suffering from <sup>high</sup> bias, getting more training data will not (by itself) help much

### III Machine Learning System Design | Prioritizing what to work on

#### Building a spam classifier

- Supervised learning  $x = \text{features of email}$ ,  $y = \text{spam}(1)$  or  $\text{spam}(0)$ . Features  $x$ : Choose 100 words indicative of spam/not spam

Note: In practice, take most frequently occurring  $n$  words (10,000 to 50,000) in training set, rather than manually pick 100 words

Slow to spend you time to make it have low error?

- Collect lots of data
  - Eg "honeypot" project
- Develop sophisticated features for message body, eg should "discount" and "discounts" be treated as the same word? How about "deal" and "dealer"? Features about punctuation?
- Develop sophisticated algorithm to detect misspelling (e.g. mortgage, medicine, w4ches).

#### 101.2 Error analysis

$m_{cv} = 500$  examples in cross validation set

Algorithm misclassifies 100 emails

Manually examine the 100 errors, and categorize them based on:

i) What type of email it is

ii) What cues (features) you think would have helped the algorithm classify them correctly.

Pharma

Replica/fake

Steal passwords

Other

→ Deliberate misspellings:

(mOrgage, medicine, etc)

→ Unusual email routing

→ Unusual (spamming) punctuation

The importance of numerical evaluation

Should discount / discounts / discounted / discounting be treated as the same word?

Can use "streaming" software (E.g. "Porter stemmer")  
universe / university

Error analysis may not be helpful for deciding if this is likely to improve performance. Only solution is to try it and see if it works

Need numerical evaluation (e.g. cross validation error) of algorithm's performance with and without stemming.

without stemming:      with stemming:

Distinguish upper vs. lower case (Moni / moni)

### 11.3. Error Metrics for Skewed Classes

Cancer classification example

Train logistic regression model  $h_0(x)$ . ( $y=1$  if cancer,  $y=0$  otherwise)

Find that you got 1% error on test set.

(99% correct diagnoses)

Only 0.50% of patients have cancer  
skewed classes

Precision / Recall

$y=1$  in presence of rare class that we want to detect

Actual Values			
		1	0
Predicted class	1	True Positive	False Positives
	0	False negatives	True negative

Precision

(Of all patient where we predicted  $y=1$ , what fraction has cancer?)

$$\frac{\text{True positives}}{\text{# predicted positives}} = \frac{\text{True positive}}{\text{True pos} + \text{Fake pos}}$$

Recall

(of all patients that actually have cancer who fraction did we correctly detect as having cancer)

$$\frac{\text{True positives}}{\text{# actual positive}} = \frac{\text{True positive}}{\text{True pos} + \text{Fake neg}}$$

## 11.4 Trade off precision and recall

Logistic regression:  $0 \leq h_0(x) \leq 1$

Predict 1 if  $h_0(x) \geq 0.5$

Predict 0 if  $h_0(x) < 0.5$

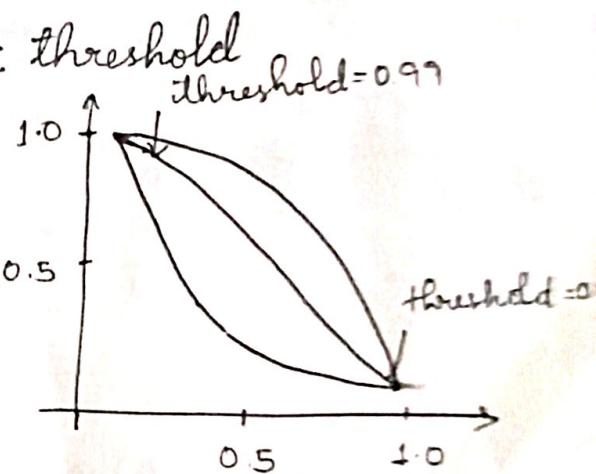
Suppose we want to predict  $y=1$  (cancer) only if very confident

→ higher precision, lower recall

Suppose we want to avoid missing too many cases of cancer (avoid false negatives).

→ higher recall, lower precision

More generally. Predict 1 if  $h_0(x) \geq \text{threshold}$



## F Score (F score)

How to compare precision/recall numbers?

	Precision (P)	Recall (R)
Algorithm 1	0.5	0.4
Algorithm 2	0.7	0.1
Algorithm 3	0.02	1.0

$$\text{Average } \frac{P+R}{2} \times$$

$$F_1 \text{ score } \frac{2 \frac{PR}{P+R}}{P+R}$$

## 11.5. System Design

Designing a high accuracy learning system

E.g. Classify between confusable words  
{to, too, too} {then, than}

→ for breakfast I ate — eggs.

## Algorithms

- Perceptron (Logistic regression)
- Winnow
- Memory-based
- Naive Bayes

## Large Data Rationale

Assume features  $x \in \mathbb{R}^{n+1}$  has sufficiently information to predict  $y$  accurately

Example: For a breakfast I ate    eggs

Counterexample: Predict housing price from only size ( $\text{feet}^2$ ) and no. other features

Useful test: Given the input  $x$ , can a human expert confidently predict  $y$ ?

→ Use a learning algorithm with many parameters (e.g. logistic regression / linear regression with many features; neural network with many hidden units).

low bias algorithm

→  $J_{\text{train}}(\theta)$  will be small

Use a very large training set (unlikely to overfit) low variance ←

→  $J_{\text{train}}(\theta) \approx J_{\text{test}}(\theta)$

$J_{\text{test}}(\theta)$  will be small